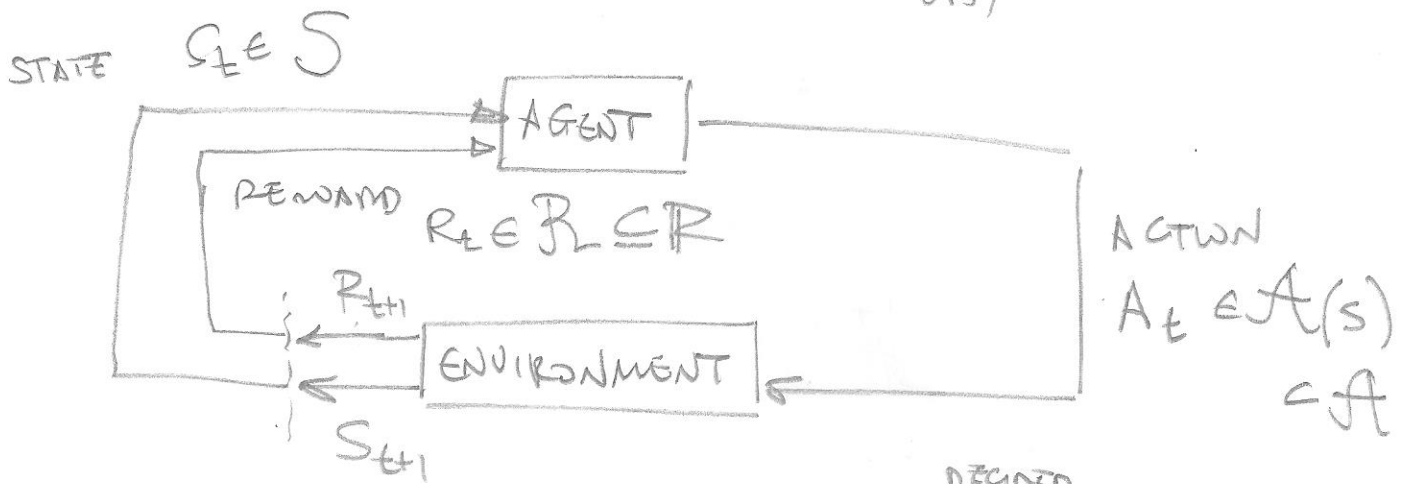


# CHAPTER 3: FINITE MARKOV DECISION PROCESSES (MDPs)

(11)



DECIDED JOINTLY, SO HAVE COMMON INDEX

TRAJECTORY:  $S_0, A_0, R_1, S_1, A_1, R_2, S_2, R_3, \dots$

MDP IS FINITE, I.E.  $|S|, |R|, |\mathcal{A}(s)| < \infty$

DYNAMIC FUNCTION

$$P(s', r | s, a) \triangleq P_r [S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a]$$

$$P: S \times R \times S \times \mathcal{A} \rightarrow [0, 1]$$

• IT IS A PMF, THEREFORE:

$$\sum_{s' \in S} \sum_{r \in R} P(s', r | s, a) = 1 \quad \forall s \in S, a \in \mathcal{A}(s)$$

• THE MDP IS HOMOGENEOUS: LEFT HAND SIDE DOES NOT DEPEND ON TIME  $t$

• THE SYSTEM HAS THE MARKOV PROPERTY:

STATISTICS OF  $S_t, R_t$  DEPEND ONLY ON  $S_{t-1}$  AND  $A_{t-1}$ . SO  $S_{t-1}$  MIGHT HAVE TO BE LARGE

STATE TRANSITION PROBABILITIES:

$$P(s' | s, a) \triangleq P_r [S_t = s' | S_{t-1} = s, A_{t-1} = a]$$

$$= \sum_{r \in R} P(s', r | s, a)$$

EXPECTED REWARDS

$$V(s, a) \triangleq E [R_t | S_{t-1} = s, A_{t-1} = a] =$$

$$\sum_{r \in R} \sum_{s' \in S} r P(s', r | s, a) =$$

$$\sum_{r \in R} r \left( \sum_{s' \in S} P(s', r | s, a) \right)$$

PROBABILITY MIT GET REWARD r

EXPECTED REWARDS FOR (STATE, ACTION, NEXT STATE) TRIPLE:

$$V(s, a, s') \triangleq E [R_t | S_{t-1} = s, A_{t-1} = a, S_t = s']$$

$$= \sum_{r \in R} r P(r | s, a, s') =$$

$$\sum_{r \in R} r \frac{P(r, s, a, s')}{P(s, a, s')} = \sum_{r \in R} r \frac{P(r, s, a, s') / P(s, a)}{P(s, a, s') / P(s, a)}$$

→

$$v(s, a, s') = \sum_{r \in R} r \frac{P(s', r | s, a)}{P(s' | s, a)}$$

(SO REWARD STATISTICS DEPEND ON NEXT STATE)

NOTE: WE USED BASIC PROPERTY

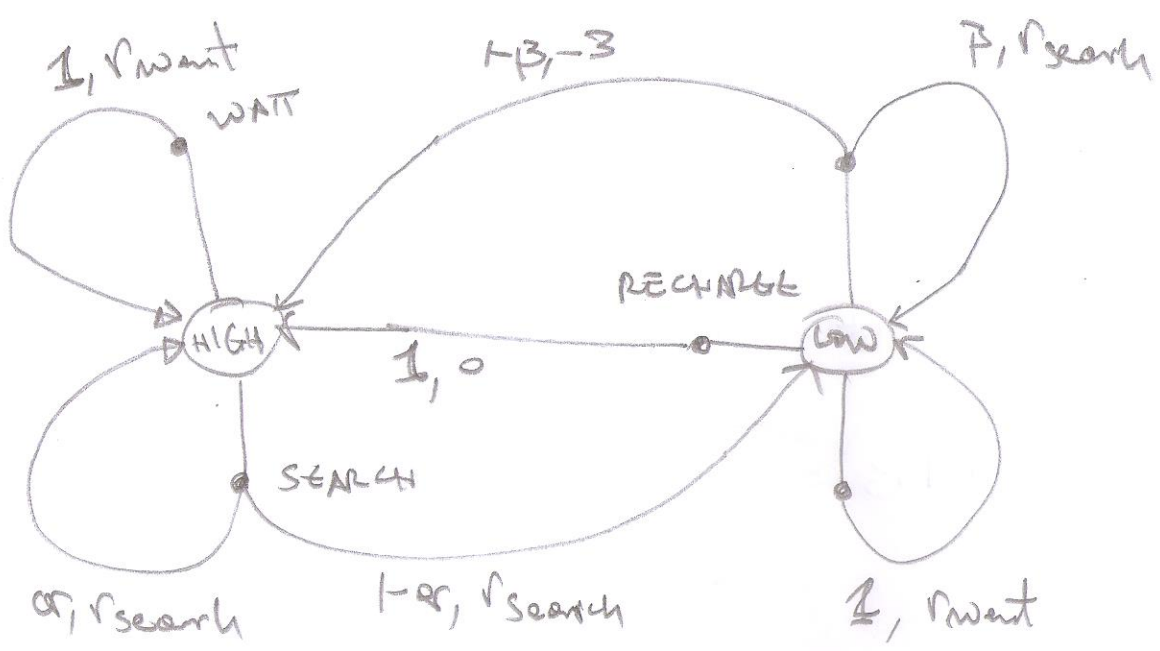
$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

EXAMPLE 3.3. A RECYCLING ROBOT

$S = \{ \text{high, low} \}$  (BATTERY LEVELS)

$A = \{ \text{search, wait, recharge} \}$

s	a	s'	P(s' s,a)	v(s,a,s')
high	search	high	$\alpha$	$r_{\text{search}}$
-  -	-  -	low	$1-\alpha$	-  -
low	search	high	$1-\beta$	-3
-  -	-  -	low	$\beta$	$r_{\text{search}}$
high	wait	high	1	$r_{\text{wait}}$
high	wait	low	0	-
low	wait	high	0	-
low	wait	low	1	$r_{\text{wait}}$
low	recharge	high	1	0
low	recharge	low	0	-



OBSERVE THAT: 1) MARKOVIAN PROPERTY IS NOT JUSTIFIED WELL

- 2)  $r_{search} > r_{want}$ , OTHERWISE SOLUTION IS TRIVIAL
- 3) WHAT IS OPTIMAL SOLUTION;

### 3.3 RETURNS AND EPISODES

WHAT DO WE WANT TO MAXIMIZE?

TWO CASES:

1) THERE IS A SPECIAL TERMINAL STATE, WHICH WE ACHIEVE AFTER FINITE R.V. T. THEN AT TIME t WE WANT TO MAXIMIZE THE EXPECTED RETURN:

$$G_t \triangleq R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$$

RETURN, IT IS AN R.V. THE EVOLUTION OF THE STATES HAPPENS IN EPISODIC TASKS, OR EPISODES (EXAMPLES: GAMES & BACK GAMMON)

② CONTINUING TASKS : THERE IS NO TERMINATING STATE AND WE GO FOREVER. AT STEP  $t$ , WE WANT TO MAXIMIZE THE EXPECTED DISCOUNTED RETURN :

$$G_t^A = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

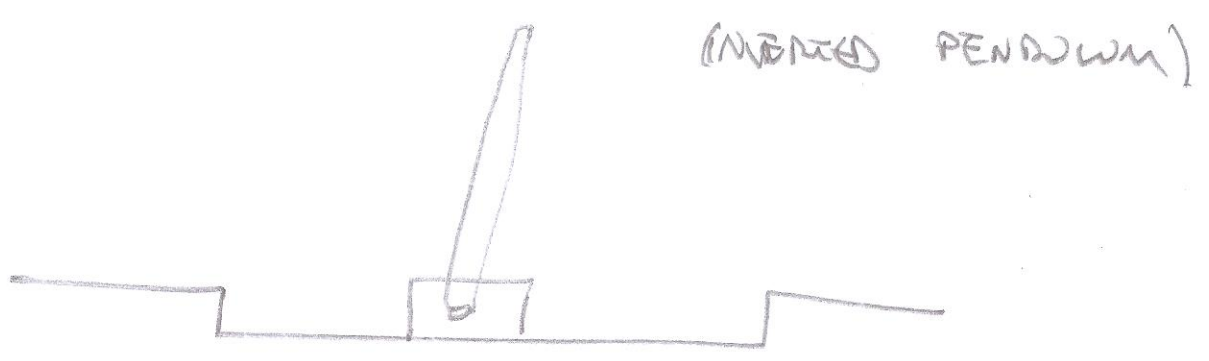
WHERE  $\gamma \in [0, 1]$  IS THE DISCOUNT RATE

A USEFUL FORMULA FOR THIS CASE :

$$\begin{aligned}
G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\
&= R_{t+1} + \gamma [R_{t+2} + \gamma R_{t+3} + \dots] \\
&= R_{t+1} + \gamma G_{t+1} \Rightarrow \boxed{G_t = R_{t+1} + \gamma G_{t+1}}
\end{aligned}$$

EXAMPLES: CHEMICAL PROCESSES WE DO NOT WANT TO END.

### EXAMPLE 3.4 POLE - BALANCING



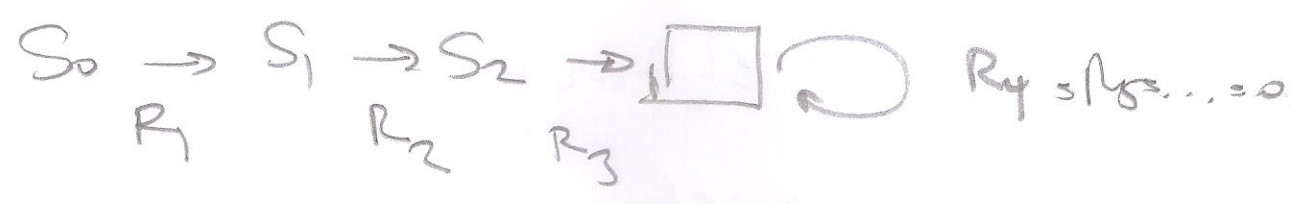
THIS CAN BE MODELED BOTH USING EPISODIC AND CONTINUING TASKS

### 3.4 UNIFIED NOTATION OF $G_t$

WE HAVE GIVEN TWO FORMULAS FOR  $G_t$ , ONE FOR EPISODIC, ONE FOR CONTINUING TASKS. WE NEED TO HAVE ONLY ONE. WE INTRODUCE TWO CONVENTIONS.

1) WE DO NOT INDEX EPISODES SEPARATELY, SO WE ONLY WRITE  $A_1, A_2, A_3, \dots$  AND NOT  $A_{1,i}, A_{2,i}, A_{3,i}, \dots$  FOR EPISODE  $i$ .

2) WE INTRODUCE ABSORBING STATE



SO, FROM NOW ON, ALWAYS,

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

BT MAYBE WE NEED  $\gamma < 1$ .

### 3.5 POLICIES AND VALUE FUNCTIONS

IF AGENT FOLLOWS POLICY  $\pi$  THEN PROBABILITY THAT  $A_t = a$  IF  $S_t = s$  IS  $\pi(a|s)$

SO THE POLICY IS WHAT WE WANT TO OPTIMIZE

OBSERVE THAT

$$E[R_{t+1} | S_t = s] = E \left[ E[R_{t+1} | S_t = s, A_t] \right]$$

$$= \sum_{\alpha} \pi(\alpha|s) v(s, \alpha)$$

$$= \sum_{\alpha} \pi(\alpha|s) \cdot \sum_{s', r} v p(s', r|s, \alpha)$$

STATE-VALUE FUNCTION FOR POLICY  $\pi$ :

$$V_{\pi}(s) \triangleq E_{\pi} [G_t | S_t = s] = E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k} | S_t = s \right]$$

$\forall s \in S$

ACTION-VALUE FUNCTION FOR POLICY  $\pi$ :

$$q_{\pi}(s, \alpha) \triangleq E_{\pi} [G_t | S_t = s, A_t = \alpha]$$

$$= E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k} | S_t = s, A_t = \alpha \right]$$

OBSERVE THAT THE TWO ARE CONNECTED:

$$V_{\pi}(s) = \sum_{\alpha} \pi(\alpha|s) q_{\pi}(s, \alpha) \quad (1)$$

$$q_{\pi}(s, \alpha) = \sum_{s', r} p(s', r|s, \alpha) [r + \gamma V_{\pi}(s')] \quad (2)$$

$$= \sum_{s', r} p(s', r|s, \alpha) E_{\pi} [G_t | S_t = s, A_t = \alpha, S_{t+1} = s', R_{t+1} = r]$$

FURTHERMORE, THEY ARE CONNECTED TO THEMSELVES,

IN SOME WAY:

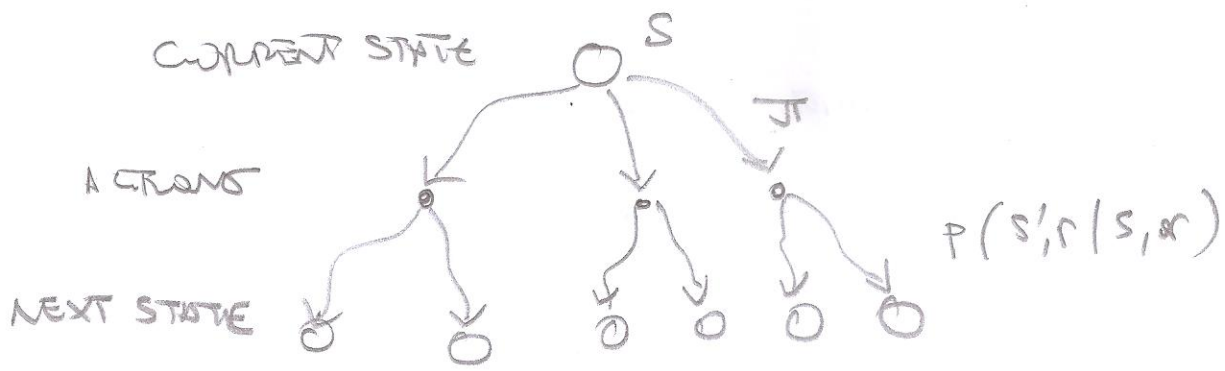
$$\begin{aligned}
V_{\pi}(s) &\triangleq E_{\pi} [G_t | S_t = s] \\
&= E_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s] \quad (3.14) \\
&= \sum_{\alpha} \pi(\alpha | s) \sum_{s'} \sum_r p(s', r | s, \alpha) \left[ r + \gamma E_{\pi} [G_{t+1} | S_{t+1} = s'] \right] \\
&= \sum_{\alpha} \pi(\alpha | s) \sum_{s'} \sum_r p(s', r | s, \alpha) \left[ r + \gamma V_{\pi}(s') \right] \quad \text{Use } \} \\
&\quad \text{(ONE OF THEM ANYWAY)}
\end{aligned}$$

$S_t = s$   
 $R_{t+1} = r$   
 $A_t = \alpha$

THIS IS THE BELLMAN EQUATION FOR  $V_{\pi}(s)$ .

IT FORMS A SYSTEM OF LINEAR EQUATIONS THAT HAS A UNIQUE SOLUTION (USUALLY, THERE IS TRICKY)

WE CAN CONCEPTUALIZE THIS USING A BACKUP DIAGRAM



THERE IS ALSO A BELLMAN EQUATION FOR

$Q_{\pi}(s, \alpha)$ .  $\sim$  SIMILAR ①, ② :



$$q_{\pi}(s, a) = \sum_{s', r} P(s', r | s, a) \left[ r + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a') \right]$$

$$= \sum_{s', r} P(s', r | s, a) \left[ r + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a') \right]$$

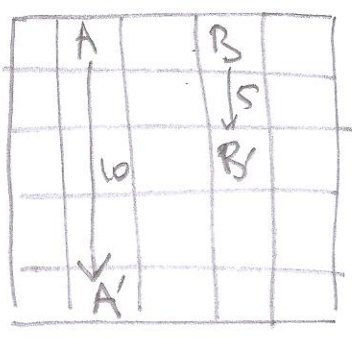
(THIS IS HARDER TO USE)

SO, HOW DO WE FIND  $v_{\pi}(s)$ ?

- 1) WE CAN DO SIMULATIONS (MONTE CARLO)
- 2) WE CAN SOLVE A LINEAR SYSTEM

OF COURSE, AIM IS TO FIND OPTIMAL POLICY, WE EXAMINE THIS NEXT.

EXAMPLE 3.5 GRIDWORLD



- 1) ACTIONS ARE UNIFORM RANDOM WALK EXCEPT FOR STATES A, B
- 2) REWARDS ARE 0 EXCEPT THOSE SHOWN ON LEFT AND THOSE TAKING AGENT OFF THE GRID, WHICH ARE -1
- 3)  $\gamma = 0.9$

THEN  $v_{\pi}(s)$  IS

OBSERVE THAT  $3.3 < 10$

$5.3 > 5$

WHY?

3.3	8.9	4.7	5.3	1.5
1.5	3.0	2.9	1.9	0.5
0.1	0.7	0.7	0.4	-2.4
-1.0	-2.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

# NESTED EXPECTATION

20

LET  $X, Y$  TWO R.V.'S WITH PMF'S

$$P_X(x), P_Y(y), P_{XY}(x,y), P_{X|Y}(x|y) = \frac{P_{XY}(x,y)}{P_Y(y)}$$

WE HAVE

$$E[X] = \sum_x x P_X(x), \quad E[X|Y=y] = \sum_x x P_{X|Y}(x|y)$$

OBSERVE THAT  $E[X|Y]$  IS ALSO AN R.V.  
THEREFORE, IT HAS AN EXPECTED VALUE:

$$E[E[X|Y]] = \sum_y P_Y(y) \sum_x x P_{X|Y}(x|y) = \textcircled{*}$$

(NOTE: •  $E[g(Y)] = \sum_y g(y) P_Y(y)$ )

•  $E[g(X,Y)] = \sum_{x,y} g(x,y) P_{XY}(x,y)$ )

$$\begin{aligned} \textcircled{*} &= \sum_{x,y} x P_Y(y) P_{X|Y}(x|y) = \sum_{x,y} x P_{XY}(x,y) \\ &= E(X) \Rightarrow \end{aligned}$$

$$E[X] = E[E[X|Y]]$$

NESTED EXPECTATION FORMULA

NICE APPLICATION: MINER'S PROBLEM